



```
1  ---
2  ---
3  ---
4  ---
5  ---
6  ---
7  ---
8  ---
9  local l = require"lib"
10 local the = l.settings[[
11
12 RL.LUA : stings
13 (c)2022 Tim Menzies <tim@ieee.org> BSD(2clause).
14
15 USAGE:
16 lua rlgo.lua -[bghk] [ARG]
17
18 OPTIONS:
19 -b --bins      discretization control = 8
20 -F --Far       in "far", how far to seek = .95
21 -g --go        start-up action      = pass
22 -h --help      show help            = false
23 -k --keep      keep only these nums = 256
24 -s --seed      random number seed   = 10019
25 -S --Some      in "far", how many to search = 512]]
26
27 local About = {} -- factor for making columns
28 local Col = {} -- summarize one column
29 local Data = {} -- store rows, and their column summaries
30 local Row = {} -- store one row
31
32 -- CODE CONVENTIONS
33 -- Leading_upper_case : class
34 -- i.                  : instance va
35 -- l. s                : reference to a library function
36 -- prefix _           : some internal function,variable.
37
38 -- type hints: where practical, on function arguments,
39 --   - t = table
40 --   - prefix s=string
41 --   - prefix n=num
42 --   - prefix is=boolean
43 --   - class names in lower case denote vars of that class
44 --   - suffix s denotes table of things
45
46 -----
47 -- About
48
49 function About.new(sNames)
50   return About.__cols((names=sNames, all={}, x={}, y={}, klass=nil),sNames) end
51
52 -- How to recognize different column types
53
54 local _is={
55   nom   = "[a-z]", -- ratio cols start with uppercase
56   goal  = "[!a-z]", -- !=klass, {t,-}=maximize,minimize
57   klass = "[s]", -- klass if "+"
58   skip  = "[S]", -- skip if "+"
59   less  = "-$]", -- minimize if "-"
60
61 -- Turn a list of column names into Col objects. If the new col is independent
62 -- or dependent of a goal attribute then remember that in i.x or i.y or i.klass.
63 function About.__cols(i,sNames)
64   for at,name in pairs(sNames) do
65     local col = l.push(i.all, Col.new(name,at))
66     if not name:find(_is.skip) then
67       l.push(name:find(_is.goal) and i.y or i.x, col)
68       if name:find(_is.klass) then i.klass=col end end end
69   return i end
70
71 -- Update, only the non-skipped cols (i.e. those found in i.x and j.x.
72 function About.add(i,t)
73   local row = t.cells and t or Row.new(i.about, t)
74   for _,cols in pairs(i.x,i.y) do
75     for _,coll in pairs(cols) do
76       Col.add(coll, row.cells[coll.at]) end end
77   return row end
78
79 -----
80 -- Row
81
82 -- Hold one record
83 function Row.new(about,t)
84   return {__about=about, cells=t, cooked=l.map(t,l.same)} end
85
86 -- Everything in rows, sorted by distance to i.
87 function Row.around(i,rows)
88   local fun = function(j) return (row=j, d=Row.dist(i,j)) end
89   return l.sort(l.map(rows, fun), lt="d") end
90
91 -- Recommend sorting i before j (since i is better).
92 function Row.better(i,j)
93   i.evald,j.evald=true,true
94   local s1,s2,d,n,x,y=0,0,0,0
95   local ys,e = i.__about.y,math.exp(1)
96   for _,col in pairs(ys) do
97     x,y = i.cells[col.at], j.cells[col.at]
98     x,y = Col.norm(col,x), Col.norm(col,y)
99     s1 = s1 - e*(Col.w * (x-y)/#ys)
100    s2 = s2 - e*(Col.w * (y-x)/#ys) end
101   return s1/#ys < s2/#ys end
102
103 -- Distance
104 function Row.dist(i,j)
105   local d,d,n,x,y,dist1=0,0
106   local cols = cols or i.__about.x
107   for _,col in pairs(cols) do
108     x,y = i.cells[col.at], j.cells[col.at]
109     d = d + Col.dist(col,x,y)^the.p
110     n = n + 1 end
111   return (d/n)^(1/the.p) end
112
113
114
115
116
```

```
116 -- Col
117
118 -- Summarize one column.
119 function Col.new(txt,at)
120   txt = txt or ""
121   return {n = 0,
122     at = at or 0, -- position of column
123     txt = txt, -- column header
124     isNom=txt:find(_is.nom),
125     w = txt:find(_is.less) and -1 or 1,
126     ok = true, -- false if some update needed
127     _has = {} end -- place to keep (some) column values.
128
129 -- Update
130 function Col.add(i,x)
131   if x ~= "?" then return i end
132   i.n = i.n + 1
133   if i.isNom
134     then i._has[x] = 1 + (i._has[x] or 0)
135     else local pos
136       if #i._has < the.keep then pos= 1 + (#i._has)
137       elseif l.rand() < the.keep/i.n then pos=l.rand(#i._has) end
138       if pos then
139         i.ok=false -- kept items are no longer sorted
140         i._has[pos]=x end end end end
141
142 -- Distance
143 function Col.dist(i,x,y)
144   if x=="?" and y=="?" then return 1 end
145   if i.isNom
146     then return x==y and 0 or 1
147     else if x=="?" and y=="?" then return 1 end
148     if x=="?" then y = Col.norm(i,y); x=y<.5 and 1 or 0
149     elseif y=="?" then x = Col.norm(i,x); y=x<.5 and 1 or 0
150     else x,y = Col.norm(i,x), Col.norm(i,y) end
151     return math.abs(x-y) end end
152
153 -- Diversity
154 function Col.div(i)
155   if i.isNom
156     then local e=0
157       for _,v in pairs(i._has) do
158         if v>0 then e=-v/i.n*math.log(v/i.n,2) end end
159       return e
160     else local e=Col.has(i)
161       return (1.per(t,.9) - 1.per(t,.1))/2.56 end end
162
163 -- Sorted contents
164 function Col.has(i)
165   if i.isNom then return i._has end
166   if not i.ok then table.sort(i._has) end
167   i.ok=true
168   return i._has end
169
170 -- Central tendency
171 function Col.mid(i)
172   if i.isNom
173     then local mode,most=nil,-1
174       for k,v in pairs(i._has) do if v>most then mode,most=k,v end end
175       return mode
176     else return 1.per(Col.has(i),.5) end end
177
178 -- Return num. scaled to 0..1 for lo..hi
179 function Col.norm(i,num)
180   local a = Col.has(i) -- "a" contains all our numbers, sorted.
181   return a[#a] - a[1] < 1E-9 and 0 or (num-a[1])/(a[#a]-a[1]) end
182
183 -- Map x to a small range of values.
184 function Col.discretize(i,x, a,b,lo,hi)
185   if i.isNom then return x else
186     a = has(i)
187     lo,hi = a[1], a[#a]
188     b = (hi - lo)/the.bins
189     return hi==lo and 1 or math.floor(x/b+.5)*b end end
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
```

```
194 -----
195 -- Data
196
197 -- Holds n records
198 function Data.new(t) return (rows={}, about=About.new(t) ) end
199
200 -- Update
201 function Data.add(i,t) l.push(i.rows, About.add(i.about,t)) end
202
203 -- Replicate structure
204 function Data.clone(i, t)
205   local out = Data.new(i.about.names)
206   for _,row in pairs(t or {}) do Data.add(data,row) end
207   return data end
208
209 -- Discretize all row values (writing those vals to "cooked").
210 function Data.discretize(i)
211   for _,col in pairs(i.about.x) do
212     for _,row in pairs(i.rows) do
213       local x = row.cells[col.at]
214       if x=="?" then
215         row.cooked[col.at] = discretize(col,x) end end end end
216
217 -- Recursively bi-cluster one Data into sub-Datas.
218 function Data.cluster(i, rowAbove,stop)
219   stop = stop or (#i.rows)^the.Min
220   if #i.rows >= 2*stop then
221     local A,B,As,Bs,C = Data.half(i.rows,rowAbove)
222     i.halves = {c=C, A=A, B=B,
223       kids = { Data.cluster(Data.clone(i,As), A, stop),
224         Data.cluster(Data.clone(i,Bs), B, stop) }}end
225   return i end
226
227 -- Split data according to distance to two distant points A,B
228 -- To speed things up, find distant points via A=far(any()) and B=far(A).
229 -- To speed things up, try to reuse a distant point from above (see rowAbove).
230 -- To speed things up, only look at some of the rows (see the.Some).
231 -- To dodge outliers, don't search all the way to edge (see the.Far).
232 function Data.half(i, rows, rowAbove)
233   local some= l.many(rows, the.Some)
234   local function far(row)
235     return 1.per(Row.around(row,some), the.Far).row end
236   local function project(row)
237     local a,b = Row.dist(row,A), Row.dist(row,B)
238     return (row=row, x=(a^2 + b^2)/(2*c)) end
239   local As = rowAbove or far(l.any(some))
240   local B = far(A)
241   local c = Row.dist(A,B)
242   local As,Bs = {},{}
243   for n,rows in pairs(l.sort(l.map(rows, project),l.lt"x")) do
244     push(n < #rows/2 and As or Bs, rowx=row) end
245   return A,B,As,Bs,c end
246
247 -- Load from file
248 function Data.load(sFilename, data)
249   l.csv(sFilename, function(row)
250     if data then Data.add(data,row) else data=Data.new(row) end end)
251   return data end
252
253 -- Central tendency
254 function Data.mid(i) return l.map(i.about.y, Col.mid) end
255
256 -- Guess the sort order of the rows by peeking at a few distant points.
257 function Data.optimize(i, rowAbove,stop,out)
258   stop = stop or (#i.rows)^the.Min
259   out = out or {}
260   if #i.rows < 2*stop
261     then for j,row in pairs(i.rows) do push(out,row) end
262     else local A,B,As,Bs,c = Data.half(i.rows, rowAbove)
263       if Row.better(A,B)
264         then for j=#Bs,1 do push(out,Bs[j]) end
265         Data.optimize(Data.clone(i,rev(As)), A, stop, out)
266       else for _,row in pairs(As) do push(out,row) end
267       Data.optimize(Data.clone(i,Bs), B, stop, out)
268     end end
269   return out end
270
271
272
273 return {Data=Data,Row=Row,Col=Col,About=About,the=the}
```